

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE
APPLICATION FOR LETTERS PATENT

**Methods and Systems for Animating Facial Features,
and Methods and Systems for Expression
Transformation**

Inventor(s):

Stephen Marschner
Brian Guenter
Sashi Raghupathy
Kirk Olynyk
Sing Bing Kang

1 **TECHNICAL FIELD**

2 This invention relates to methods and systems for modeling and rendering
3 for realistic facial animation. In particular, the invention concerns methods and
4 systems for facial image processing.

5
6 **BACKGROUND**

7 The field of computer graphics involves rendering various objects so that
8 the objects can be displayed on a computer display for a user. For example,
9 computer games typically involve computer graphics applications that generate
10 and render computer objects for display on a computer monitor or television.
11 Modeling and rendering realistic images is a continuing challenge for those in the
12 computer graphics field. One particularly challenging area within the computer
13 graphics field pertains to the rendering of realistic facial images. As an example, a
14 particular computer graphics application may render a display of an individual
15 engaging in a conversation. Often times, the ultimately rendered image of this
16 individual is very obviously a computer-rendered image that greatly differs from a
17 real individual.

18 Modeling and rendering realistic faces and facial expressions is a
19 particularly difficult task for two primary reasons. First, the human skin has
20 reflectance properties that are not well modeled by the various shading models that
21 are available for use. For example, the well-known Phong model does not model
22 human skin very well. Second, when rendering facial expressions, the slightest
23 deviation from what would be perceived as "real" facial movement is perceived by
24 even the casual viewer as being incorrect. While current facial motion capture
25 systems can be used to create quite convincing facial animation, the captured

1 motion is much less convincing, and frequently very strange, when applied to
2 another face. For example, if a person provides a sampling of their facial
3 movements, then animating their specific facial movements is not difficult
4 considering that the face from which the movements originated is the same face.
5 Because of this, there will be movement characteristics that are the same or very
6 similar between expressions. Translating this person's facial movements to
7 another person's face, however, is not often times convincing because of, among
8 other things, the inherent differences between the two faces (e.g. size and shape of
9 the face).

10 Accordingly, this invention arose out of concerns associated with providing
11 improved systems and methods for modeling texture and reflectance of human
12 skin. The invention also arose out of concerns associated with providing systems
13 and methods for reusing facial motion capture data by transforming one person's
14 facial motions into another person's facial motions.

15

16 SUMMARY

17 The illustrated and described embodiments propose inventive techniques
18 for capturing data that describes 3-dimensional (3-D) aspects of a face,
19 transforming facial motion from one individual to another in a realistic manner,
20 and modeling skin reflectance.

21 In the described embodiment, a human subject is provided and multiple
22 different light sources are utilized to illuminate the subject's face. One of the light
23 sources is a structured light source that projects a pattern onto the subject's face.
24 This structured light source enables one or more cameras to capture data that
25 describes 3-D aspects of the subject's face. Another light source is provided and is

1 used to illuminate the subject's face. This other light source is sufficient to enable
2 various reflectance properties of the subject's face to be ascertained. The other
3 light source is used in conjunction with polarizing filters so that the specular
4 component of the face's reflectance is eliminated, i.e. only the diffuse component
5 is captured by the camera. The use of the multiple different light sources enables
6 both structure and reflectance properties of a face to be ascertained at the same
7 time. By selecting the light sources carefully, for example, by making the light
8 sources narrowband and using matching narrowband filters on the cameras, the
9 influence of ambient sources of illumination can be eliminated.

10 Out of the described illumination process, two useful items are produced—
11 (1) a range map (or depth map) and (2) an image of the face that does not have the
12 structured light source pattern in it. A 3D surface is derived from the range map
13 and surface normals to the 3D surface are computed. The processing of the range
14 map to define the 3D surface can optionally include a filtering step in which a
15 generic face template is combined with the range map to reject undesirable noise.
16 The computed surface normals and the image of the face are then used to derive an
17 albedo map. An albedo map is a special type of texture map in which each sample
18 describes the diffuse reflectance of the surface of a face at a particular point on the
19 surface. Accordingly, at this point in the process, information has been ascertained
20 that describes the 3D-aspects of a face (i.e. the surface normals), and information
21 that describes the face's reflectance (i.e. the albedo map).

22 In one embodiment, the information or data that was produced in the
23 illumination process is used to transform facial expressions of one person into
24 facial expressions of another person. In this embodiment, the notion of a code
25 book is introduced and used.

1 A code book contains data that describes many generic expressions of
2 another person (person A). One goal is to take the code book expressions and use
3 them to transform the expressions of another person (person B). To do this, an
4 inventive method uses person B to make a set of training expressions. The
5 training expressions consist of a set of expressions that are present in the code
6 book. By using the training expressions and each expression's corresponding code
7 book expression, a transformation function is derived. The transformation
8 function is then used to derive a set of synthetic expressions that should match the
9 expressions of person B. That is, once the transformation function is derived, it is
10 applied to each of the expressions in the code book so that the code book
11 expressions match the expressions of person B. Hence, when a new expression is
12 received, e.g. from person B, that might not be in the training set, the synthesized
13 code book expressions can be searched for an expression that best matches the
14 expression of person B.

15 In another embodiment, a common face structure is defined that can be
16 used to transform facial expressions and motion from one face to another. In the
17 described embodiment, the common face structure comprises a coarse mesh
18 structure or "base mesh" that defines a subdivision surface that is used as the basis
19 for transforming the expressions of one person into another. A common base
20 mesh is used for all faces thereby establishing a correspondence between two or
21 more faces. Accordingly, this defines a structure that can be used to adapt face
22 movements from one person to another. According to this embodiment, a
23 technique is used to adapt the subdivision surface to the face model of a subject.
24 The inventive technique involves defining certain points on the subdivision
25 surface that are mapped directly to corresponding points on the face model. This

1 is true for every possible different face model. By adding this constraint, the base
2 mesh has a property in that it fits different face models in the same way. In
3 addition, the inventive algorithm utilizes a smoothing functional that is minimized
4 to ensure that there is a good correspondence between the base mesh and the face
5 model.

6 In another embodiment, a reflectance processing technique is provided that
7 gives a measure of the reflectance of the surface of a subject's face. To measure
8 reflectance, the inventive technique separates the reflectance into its diffuse and
9 specular components and focuses on the treatment of the diffuse components.

10 To measure the diffuse component, an albedo map is first defined. The
11 albedo map is defined by first providing a camera and a subject that is illuminated
12 by multiple different light sources. The light sources are filtered by polarizing
13 filters that, in combination with a polarizing filter placed in front of the camera,
14 suppress specular reflection or prevent specular reflection from being recorded. A
15 sequence of images is taken around the subject's head. Each individual image is
16 processed to provide an individual albedo map that corresponds to that image. All
17 of the albedo maps for a particular subject are then combined to provide a single
18 albedo map for the subject's entire face.

19

20 **BRIEF DESCRIPTION OF THE DRAWINGS**

21 Fig. 1 is a high level diagram of a general purpose computer that is suitable
22 for use in implementing the described embodiments.

23 Fig. 2 is a schematic diagram of a system that can be utilized to capture
24 both structural information and reflectance information of a subject's face at the
25 same time.

1 Fig. 3 is a flow diagram that describes an exemplary method for capturing
2 structural information and reflectance information in accordance with the
3 described embodiment.

4 Fig. 4 is a schematic diagram that illustrates an exemplary code book and
5 transformation function in accordance with the described embodiment.

6 Fig. 5 is a flow diagram that illustrates an expression transformation
7 process in accordance with the described embodiment.

8 Fig. 6 is a high level diagram of an exemplary system in which certain
9 principles of the described embodiments can be employed.

10 Fig. 7 is a collection of exemplary color plates that illustrate an exemplary
11 expression transformation in accordance with the described embodiment.

12 Fig. 8 is a color picture that illustrates the process of mapping the same
13 subdivision control mesh to a displaced subdivision surface for different faces.

14 Fig. 9 is a color picture that illustrates exemplary constraints that are
15 utilized to enforce feature correspondence during surface fitting.

16 Fig. 10 is a flow diagram that describes steps in a surface fitting method in
17 accordance with the described embodiment.

18 Fig. 11 is a schematic diagram of an exemplary system that can be
19 employed to build an albedo map for a face in accordance with the described
20 embodiment.

21 Fig. 12 is a color picture of an exemplary albedo map for two photographs
22 that are projected into texture space and corrected for lighting.

23 Fig. 13 is a color picture of an exemplary weighting function that
24 corresponds to the Fig. 12 photographs.

25 Fig. 14 is a color picture of two full albedo maps for two different data sets.

1 Fig. 15 is a color diagram of the Fig. 14 albedo maps after editing.

2 Fig. 16 is a collection of color pictures of a face model that is rendered in
3 different orientations and under different lighting conditions.

4 Fig. 17 is a flow diagram that describes steps in a method for creating an
5 albedo map in accordance with the described embodiment.

6 Fig. 18 is a flow diagram that describes steps in a method for computing an
7 albedo for a single pixel in accordance with the described embodiment.

8

9 **DETAILED DESCRIPTION**

10 **Overview**

11 Rendering realistic faces and facial expressions requires very good models
12 for the reflectance of skin and the motion of the face. Described below are
13 methods and techniques for modeling, animating, and rendering a face using
14 measured data for geometry, motion, and reflectance that realistically reproduces
15 the appearance of a particular person's face and facial expressions. Because a
16 complete model is built that includes geometry and bi-directional reflectance, the
17 face can be rendered under any illumination and viewing conditions. The
18 described modeling systems and methods create structured face models with
19 correspondences across different faces, which provide a foundation for a variety of
20 facial animation operations.

21 The inventive embodiments discussed below touch upon each of the parts
22 of the face modeling process. To create a structured, consistent representation of
23 geometry that forms the basis for a face model and that provides a foundation for
24 many further face modeling and rendering operations, inventive aspects extend
25 previous surface fitting techniques to allow a generic face to be conformed to

1 different individual faces. To create a realistic reflectance model, the first known
2 practical use of recent skin reflectance measurements is made. In addition, newly
3 measured diffuse texture maps have been added using an improved texture capture
4 process. To animate a generic mesh, improved techniques are used to produce
5 surface shapes suitable for high quality rendering.

6

7 Exemplary Computer System

8 Preliminarily, Fig. 1 shows a general example of a desktop computer 130
9 that can be used in accordance with the described embodiments. Various numbers
10 of computers such as that shown can be used in the context of a distributed
11 computing environment. These computers can be used to render graphics and
12 process images in accordance with the description given below.

13 Computer 130 includes one or more processors or processing units 132, a
14 system memory 134, and a bus 136 that couples various system components
15 including the system memory 134 to processors 132. The bus 136 represents one
16 or more of any of several types of bus structures, including a memory bus or
17 memory controller, a peripheral bus, an accelerated graphics port, and a processor
18 or local bus using any of a variety of bus architectures. The system memory 134
19 includes read only memory (ROM) 138 and random access memory (RAM) 140.
20 A basic input/output system (BIOS) 142, containing the basic routines that help to
21 transfer information between elements within computer 130, such as during start-
22 up, is stored in ROM 138.

23 Computer 130 further includes a hard disk drive 144 for reading from and
24 writing to a hard disk (not shown), a magnetic disk drive 146 for reading from and
25 writing to a removable magnetic disk 148, and an optical disk drive 150 for

1 reading from or writing to a removable optical disk 152 such as a CD ROM or
2 other optical media. The hard disk drive 144, magnetic disk drive 146, and optical
3 disk drive 150 are connected to the bus 136 by an SCSI interface 154 or some
4 other appropriate peripheral interface. The drives and their associated computer-
5 readable media provide nonvolatile storage of computer-readable instructions, data
6 structures, program modules and other data for computer 130. Although the
7 exemplary environment described herein employs a hard disk, a removable
8 magnetic disk 148 and a removable optical disk 152, it should be appreciated by
9 those skilled in the art that other types of computer-readable media which can
10 store data that is accessible by a computer, such as magnetic cassettes, flash
11 memory cards, digital video disks, random access memories (RAMs), read only
12 memories (ROMs), and the like, may also be used in the exemplary operating
13 environment.

14 A number of program modules may be stored on the hard disk 144,
15 magnetic disk 148, optical disk 152, ROM 138, or RAM 140, including an
16 operating system 158, one or more application programs 160, other program
17 modules 162, and program data 164. A user may enter commands and
18 information into computer 130 through input devices such as a keyboard 166 and a
19 pointing device 168. Other input devices (not shown) may include a microphone,
20 joystick, game pad, satellite dish, scanner, and one or more cameras, or the like.
21 These and other input devices are connected to the processing unit 132 through an
22 interface 170 that is coupled to the bus 136. A monitor 172 or other type of
23 display device is also connected to the bus 136 via an interface, such as a video
24 adapter 174. In addition to the monitor, personal computers typically include other
25 peripheral output devices (not shown) such as speakers and printers.

1 Computer 130 commonly operates in a networked environment using
2 logical connections to one or more remote computers, such as a remote computer
3 176. The remote computer 176 may be another personal computer, a server, a
4 router, a network PC, a peer device or other common network node, and typically
5 includes many or all of the elements described above relative to computer 130,
6 although only a memory storage device 178 has been illustrated in Fig. 1. The
7 logical connections depicted in Fig. 1 include a local area network (LAN) 180 and
8 a wide area network (WAN) 182. Such networking environments are
9 commonplace in offices, enterprise-wide computer networks, intranets, and the
10 Internet.

11 When used in a LAN networking environment, computer 130 is connected
12 to the local network 180 through a network interface or adapter 184. When used
13 in a WAN networking environment, computer 130 typically includes a modem 186
14 or other means, such as a network interface, for establishing communications over
15 the wide area network 182, such as the Internet. The modem 186, which may be
16 internal or external, is connected to the bus 136 via a serial port interface 156. In a
17 networked environment, program modules depicted relative to the personal
18 computer 130, or portions thereof, may be stored in the remote memory storage
19 device. It will be appreciated that the network connections shown are exemplary
20 and other means of establishing a communications link between the computers
21 may be used.

22 Generally, the data processors of computer 130 are programmed by means
23 of instructions stored at different times in the various computer-readable storage
24 media of the computer. Programs and operating systems are typically distributed,
25 for example, on floppy disks or CD-ROMs. From there, they are installed or

1 loaded into the secondary memory of a computer. At execution, they are loaded at
2 least partially into the computer's primary electronic memory. The invention
3 described herein includes these and other various types of computer-readable
4 storage media when such media contain instructions or programs for implementing
5 the steps described below in conjunction with a microprocessor or other data
6 processor. The invention also includes the computer itself when programmed
7 according to the methods and techniques described below.

8 For purposes of illustration, programs and other executable program
9 components such as the operating system are illustrated herein as discrete blocks,
10 although it is recognized that such programs and components reside at various
11 times in different storage components of the computer, and are executed by the
12 data processor(s) of the computer.

13

14 **Exemplary System for Capturing Structure and Properties of a Facial**
Surface

15 In the past, capturing systems have not been able to capture both facial
16 structure and reflectance properties of a whole face independently at the same
17 time. There are systems that, for example, use structured light to capture the
18 structure of the face--but these systems do not capture properties of the face such
19 as the reflectance. Similarly, there are systems that capture reflectance of the face--
20 but such systems do not capture facial structure. The ability to capture facial
21 structure and reflectance independently at the same time makes it possible to
22 perform additional operations on collected data which is useful in various face
23 rendering and animation operations. One particular example of an exemplary
24 rendering operation is described below. It is to be understood, however, that the
25 information or data that is produced as a result of the system and method described

1 below can be utilized in various other areas. For example, areas of application
2 include, without limitation, recognition of faces for security, personal user
3 interaction, etc., building realistic face models for animation in games, movies,
4 etc., and allowing a user to easily capture his/her own face for use in interactive
5 entertainment or business communication.

6 Fig. 2 shows an exemplary system 200 that is suitable for use in
7 simultaneously or contemporaneously capturing facial structure and reflectance
8 properties of a subject's face. The system includes a data-capturing system in the
9 form of one or more cameras, an exemplary one of which is camera 202. Camera
10 202 can include a CCD image sensor and related circuitry for operating the array,
11 reading images from it, converting the images to digital form, and communicating
12 those images to the computer. The system also includes a facial illumination
13 system in the form of multiple light sources or projectors. In the case where
14 multiple cameras are used, they are genlocked to allow simultaneous capture in
15 time. In the illustrated example, two light sources 204, 206 are utilized. Light
16 source 204 desirably produces a structured pattern that is projected onto the
17 subject's face. Light source 204 can be positioned at any suitable location. This
18 pattern enables structural information or data pertaining to the 3-D shape of the
19 subject's face to be captured by camera 202. Any suitable light source can be
20 used, although a pattern composed of light in the infrared region can be
21 advantageously employed. Light source 206 desirably produces light that enables
22 camera 202 to capture the diffuse component of the face's reflectance property.
23 Light source 206 can be positioned at any suitable location although it has been
24 advantageously placed in line with the camera's lens 202a through, for example,
25 beam splitting techniques. This light source could also be adapted so that it

1 encircles the camera lens. This light source is selected so that the specular
2 component of the reflectance is suppressed or eliminated. In the illustrated
3 example, a linear polarizing filter is employed to produce polarized illumination,
4 and a second linear polarizer, which is oriented perpendicularly to the first, is
5 placed in front of the lens 202a so that specular reflection from the face is not
6 recorded by the camera. The above-described illumination system has been
7 simulated using light sources at different frequencies, e.g. corresponding to the red
8 and green channels of the camera. Both of the channels can, however, be in the
9 infrared region. Additionally, by selecting the light sources to be in a narrow band
10 (e.g. 780-880 nm), the influence of ambient light can be eliminated. This property
11 is only achieved when the camera is also filtered to a narrow band. Because the
12 illumination from the light source is concentrated into a narrow band of
13 wavelengths whereas the ambient light is spread over a broad range of
14 wavelengths, the light from the source will overpower the ambient light for those
15 particular wavelengths. The camera, which is filtered to record only the
16 wavelengths emitted by the source, will therefore be relatively unaffected by the
17 ambient light. As a result, the camera will only detect the influence of the selected
18 light sources on the subject.

19 Using the multiple different light sources, and in particular, an infrared light
20 source in combination with a polarized light source (which can be an infrared light
21 source as well) enables the camera (which is configured with a complementary
22 polarizer) to simultaneously or contemporaneously capture structural information
23 or data about the face (from light source 204) and reflectance information or data
24 about the face (from light source 206) independently. The structural information
25 describes 3-dimensional aspects of the face while the reflectance information

1 describes diffuse reflectance properties of the face. This information is then
2 processed by a computerized image processor, such as computer 208, to provide
3 information or data that can be used for further facial animation operations. In the
4 example about to be described, this information comprises 3-dimensional data (3D
5 data) and an albedo map.

6 Fig. 3 is a flow diagram that describes steps in a method in accordance with
7 this described embodiment. The described method enables information or data
8 that pertains to structure and reflection properties of a face to be collected and
9 processed at the same time. Step 300 illuminates a subject's face with multiple
10 different light sources. An exemplary system for implementing this step is shown
11 in Fig. 2. It will be appreciated that although two exemplary light sources are
12 utilized in the given example, other numbers of light sources can conceivably be
13 used. Step 302 measures range map data (depth map data) and image data from
14 the illumination of step 300. That is, the illumination of step 300 enables the
15 camera to detect light reflectance that is utilized to provide both range map data
16 and image data (i.e. reflectance) that does not contain the structure light source
17 pattern in it. The range map data and image data are provided to computer 208
18 (Fig. 2) for processing. At this point, step 304 can optionally apply a generic face
19 template to the range map data to reject various noise that can be associated with
20 the range map data. A generic face template can be considered as a 3D filter that
21 rejects noise in the range map data. Generic face templates will be understood by
22 those skilled in the art.

23 Step 306 uses the range map data to derive or compute a 3D surface. Any
24 suitable algorithm can be used and will be apparent to those skilled in the art.
25 Exemplary algorithms are described in the following papers: Turk & Levoy,

1 *Zippered Polygon Meshes from Range Images*, SIGGRAPH 94; F. Bernardini, J.
2 Mittleman, H. Rushmeier, C. Silva, and G. Taubin, *The Ball-Pivoting Algorithm*
3 for Surface Reconstruction, Trans. Vis. Comp. Graph. 5:4 (1999). Step 308 then
4 computes surface normal vectors (“surface normals”) to the 3D surface of step 306
5 using known algorithms. One way to accomplish this task is to compute the
6 normals to the triangles, average those triangle normals around each vertex to
7 make vertex normals, and then interpolate the vertex normals across the interior of
8 each triangle. Other methods can, of course, be utilized. Step 310 then uses the
9 computed surface normals of step 308 and the image data of step 302 to derive an
10 albedo map. An albedo is a special type of texture map in which each sample
11 describes the diffuse reflectance of the surface of a face at a particular point on the
12 facial surface. The derivation of an albedo map, given the information provided
13 above, will be understood by those skilled in the art. An exemplary algorithm is
14 described in Marschner, *Inverse Rendering for Computer Graphics*, PhD thesis,
15 Cornell University, August 1998.

16 At this point, and as shown in Fig. 2, the illumination processing has
17 produced 3D data that describes the structural features of a subject’s face and
18 albedo map data that describes the diffuse reflectance of the facial surface.

19 The above illumination processing can be used to extract the described
20 information, which can then be used for any suitable purpose. In one particularly
21 advantageous embodiment, the extracted information is utilized to extract and
22 recognize a subject’s expressions. This information can then be used for
23 expression transformation. In the inventive embodiment described just below, the
24 expressions of one person can be used to transform the expressions of another
25 person in a realistic manner.

1

2 Expression Transformation Using a Code Book

3

4 In one expression transformation embodiment, the notion of a code book is
5 introduced and is utilized in the expression transformation operation that is
6 described below. Fig. 4 shows an exemplary code book 400 that contains many
7 different expressions that have been captured from a person. These expressions
8 can be considered as generic expressions, or expressions from a generic person
9 rather than from a particular individual. In the example, the expressions range
10 from Expression 1 through Expression N. Expression 1 could be, for example, a
11 smile; Expression 2 could be a frown; Expression 3 could be an "angry"
12 expression, and the like. The expressions that are contained in code book 400 are
13 mathematically described in terms of their geometry and can be captured in any
suitable way such as the process described directly above.

14 To effect expression transformation, a transformation function is first
15 derived using some of the expressions in code book 400. To derive the
16 transformation function, the notion of a training set of expressions 402 is
17 introduced. The expression training set 402 consists of a set of expressions that
18 are provided by an individual other than the individual whose expressions are
19 described in the code book 400. The training expressions of training set 402 are a
20 subset of the code book expressions. That is, each expression in the training set
21 corresponds to an expression in the code book 400. For example, the training set
22 402 might consist of three expressions—Expression 1, Expression 2, and
23 Expression 3, where the expressions are "smile", "frown" and "angry"
24 respectively. The goal of the transformation function is to take the geometric
25 deformations that are associated with expressions of the training set, and apply

1 them to all of the expressions of the code book 400 so that the code book
2 expressions are realistic representations of the expressions. That is, consider that
3 each person's face geometrically deforms differently for any given expression. If
4 one person's geometric facial deformations for a given expression were to be
5 simply applied to another person's face for the purpose of rendering the
6 expression, the face to which the deformations were applied would likely look
7 very distorted. This is a result of not only different facial geometries, but also of
8 differing facial deformations as between the faces. Accordingly, a transformation
9 function is derived that gives the best transformation from one set of expressions
10 to another.

11 Consider again Fig. 4 where a linear transformation processor 406 is
12 shown. Transformation processor 406 can be implemented in any suitable
13 hardware, software, firmware, or combination thereof. In the illustrated example,
14 the linear transformation processor 406 is implemented in software. The linear
15 transformation processor receives as input the training set of expressions 402 and
16 the corresponding code book expressions 404. The transformation processor
17 processes the inputs to derive a transformation function 408. The transformation
18 function 408 can then be applied to all of the expressions in the code book 400 to
19 provide a synthesized set of expressions 410. The synthesized set of expressions
20 represents expressions of the code book that have been manipulated by the
21 geometric deformations associated with the expressions of the person that
22 provided the training set of expressions.

23 Facial displacements for identical expressions will not be the same on
24 different people for two reasons. First, the motion capture sample points (one
25 particular example of how one could represent face movements in this particular

1 algorithm) will not precisely correspond because of errors in placement. Second,
2 head shape and size varies from person to person.

3 The first mismatch can be overcome by resampling the motion capture
4 displacement data for all faces at a fixed set of positions on a generic mesh. This
5 is described below in more detail in the section entitled “Exemplary System and
6 Method for Building a Face Model.” There, the fixed set of positions is referred to
7 as the “standard sample positions”. The resampling function is the mesh
8 deformation function. The standard sample positions are the vertices of the face
9 mesh that correspond to the vertices of the generic mesh subdivided once.

10 The second mismatch requires transforming displacement data from one
11 face to another to compensate for changes in size and shape of the face. In the
12 illustrated example, this is done by finding a small training set of corresponding
13 expressions for the two data sets and then finding the best linear transformation
14 from one to another. As an example, consider the following: In an experimental
15 environment, emotion expressions were manually labeled for 49 corresponding
16 expressions including various intensities of several expressions. For speech
17 motion, 10,000 frames were automatically aligned using time warping techniques.

18 Each expression is represented by a $3m$ -vector g that contains all of the x , y ,
19 and z displacements at the m standard sample positions. Given a set of n
20 expression vectors for the face to be transformed, $g_{a1\dots n}$, and a corresponding set of
21 vectors for the target face, $g_{b1\dots n}$, a set of linear predictors a_j is computed, one for
22 each coordinate of g_a , by solving $3m$ linear least squares systems:

$$24 a_j \cdot g_{ai} = g_{bi}[j] \quad i = 1 \dots n$$

1 In the illustrated example, only a small subset of the points of each g_{aj} are
2 used. Specifically, those points that share edges with the standard sample point
3 under consideration. In the mesh that was used, the average valence is about 6 so
4 that the typical g_{aj} has 18 elements. The resulting system is roughly n by 18.

5 The resulting linear system may be ill-conditioned, in which case the linear
6 predictors a_j do not generalize well. The spread of the singular values is
7 controlled when computing the pseudoinverse to solve for the a_j , which greatly
8 improves generalization. All singular values less than $\alpha\sigma_1$, where σ_1 is the largest
9 singular value of the matrix and $\alpha = 0.2\dots0.1$ are zeroed out.

10 Fig. 5 is a flow diagram that describes steps in an expression transformation
11 method in accordance with this described embodiment. Step 500 provides a code
12 book of expressions. An example of such a code book is given above. Step 502
13 provides a training set of expressions. Typically, this training set is a set of
14 expressions from a person who is different from the person who provided the code
15 book expressions. The training set of expressions can be captured in any suitable
16 way. As an example, the expressions can be captured using a system such as the
17 one illustrated in Fig. 2. After the training set of expressions is provided, step 504
18 derives a transformation function using the training set and the code book. One
19 exemplary way of accomplishing this task was described above. Other methods
20 could, of course, be used without departing from the spirit and scope of the
21 claimed subject matter. For example, one could use various kinds of nonlinear
22 transformations such as neural networks, or weighted sums of basis expressions.
23 Once the transformation function is derived, it is applied to all of the expressions
24 in the code book to provide or define a synthetic set of expressions that can then
25 serve as a basis for subsequent facial animation operations.

1 2 **Exemplary Application**

3 Fig. 6 shows a system 600 that illustrates but one example of how the
4 expression transformation process described above can be employed. System 600
5 includes a transmitter computing system or transmitter 602 and a receiver
6 computing system or receiver 604 connected for communication by a network 603
7 such as the Internet. Transmitter 602 includes an illumination system 200 (Fig. 2)
8 that is configured to capture the expressions of a person as described in connection
9 with Fig. 2. Transmitter 602 also includes a code book 400, such as the one
10 described in connection with Fig. 4. It is assumed that the code book has been
11 synthesized into a synthetic set of expressions as described above. That is, using a
12 training set of expressions provided by the person whose expressions illumination
13 system 200 is configured to capture, the code book has been processed to provide
14 the synthesized set of expressions.

15 Receiver 604 includes a reconstruction module 606 that is configured to
16 reconstruct facial images from data that is received from transmitter 602.
17 Receiver 604 also includes a code book 400 that is identical to the code book that
18 is included with the transmitter 602. Assume now, that the person located at
19 transmitter 602 attempts to communicate with a person located at receiver 604. As
20 the person located at the transmitter 602 moves their face to communicate, their
21 facial expressions and movement are captured and processed by the transmitter
22 602. This processing can include capturing their expressions and searching the
23 synthesized code book to find the nearest matching expression in the code book.
24 When a matching expression is found in the synthesized code book, an index of

1 that expression can be transmitted to receiver 604 and an animated face can be
2 reconstructed using the reconstruction module 606.

3

4 **Exemplary Facial Transformation**

5 Fig. 7 shows some effects of expression transfer in accordance with the
6 described embodiment. The pictures in the first row constitute a synthetic face of
7 a first person (person A) that shows three different expressions. These pictures are
8 the result of the captured facial motion of person A. Face motion for a second
9 person (person B) was captured. The captured face motion for person B is shown
10 in the third row. Here, the 3D motion data was captured by placing a number of
11 colored dots on the person's face and measuring the dots' movements when the
12 person's face was deformed, as will be understood by those of skill in the art.
13 Motion data can, however, be captured by the systems and methods described
14 above. Person B's captured motions were then used, as described above, to
15 transform the expressions of person A. The result of this operation is shown in the
16 second row. The expressions in the three sets of pictures all correspond with one
17 another. Notice how the expressions in the first and second row look very similar
18 even though they were derived from two very different people, while the original
19 expressions of the second person (row 3) look totally unlike those of the first and
20 second rows.

21

22 **Exemplary System and Methods for Building a Face Model**

23 The model of a face that is needed to produce a realistic image has two
24 parts to it. The first part of the model relates to the geometry of the face (i.e. the
25 shape of the surface of the face) while the second part of the model relates to the

reflectance of the face (i.e. the color and reflective properties of the face). This section deals with the first part of that model—the geometry of the face.

The geometry of the face consists of a skin surface plus additional surfaces for the eyes. In the present example, the skin surface is derived from a laser range scan of the head and is represented by a subdivision surface with displacement maps. The eyes are a separate model that is aligned and merged with the skin surface to produce a complete face model suitable for high quality rendering.

Mesh Fitting

The first step in building a face model is to create a subdivision surface that closely approximates the geometry measured by the laser range scanner. In the illustrated example, the subdivision surfaces are defined from a coarse triangle mesh using Loop's subdivision rules. Loop's subdivision rules are described in detail in Charles Loop, *Smooth Subdivision Surfaces Based on Triangles*, PhD thesis, University of Utah, August 1987. In addition, the subdivision surfaces are defined with the addition of sharp edges similar to those described by Hoppe et al., *Piecewise smooth surface reconstruction*, Computer Graphics (SIGGRAPH '94 Proceedings) pps. 295-302, July 1994. Note that the non-regular crease masks are not used. In addition, when subdividing an edge between a dart and a crease vertex, only the new edge adjacent the crease vertex is marked as a sharp edge.

A single base mesh is used to define the subdivision surfaces for all of the face models, with only the vertex positions varying to adapt to the shape of each different face. In the illustrated example, a base mesh having 227 vertices and 416 triangles was defined to have the general shape of a face and to provide greater detail near the eyes and lips, where the most complex geometry and motion occur.

1 The mouth opening is a boundary of the mesh, and is kept closed during the fitting
2 process by tying together the positions of the corresponding vertices on the upper
3 and lower lips. The base mesh has a few edges marked for sharp subdivision rules
4 that serve to create corners at the two sides of the mouth opening and to provide a
5 place for the sides of the nose to fold. Because the modified subdivision rules
6 only introduce creases for chains of at least three sharp edges, this model does not
7 have creases in the surface; only isolated vertices fail to have well-defined limit
8 normals.

9 Fig. 8 shows an example of a coarse defined mesh (the center figure) that
10 was used in accordance with this example. Fig. 8 visually shows how the coarse
11 mesh can be used to map the same subdivision control (coarse) mesh to a
12 displaced subdivision surface for each face so that the result is a natural
13 correspondence from one face to another. This aspect is discussed in more detail
14 below.

15 The process used to fit the subdivision surface to each face is based on an
16 algorithm described by Hoppe et al. *Piecewise smooth surface reconstruction*,
17 Computer Graphics (SIGGRAPH '94 Proceedings) pps. 295-302, July 1994.
18 Hoppe's surface fitting method can essentially be described as consisting of three
19 phases: a topological type estimation (phase 1), a mesh optimization (phase 2), and
20 a piecewise smooth surface optimization (phase 3).

21 Phase 1 constructs a triangular mesh consisting of a relatively large number
22 of triangles given an unorganized set of points on or near some unknown surface.
23 This phase determines the topological type of the surface and produces an initial
24 estimate of geometry. Phase 2 starts with the output of phase 1 and reduces the
25 number of triangles and improves the fit to the data. The approach is to cast the

1 problem as optimization of an energy function that explicitly models the trade-off
2 between the competing goals of concise representation and good fit. The free
3 variables in the optimization procedure are the number of vertices in the mesh,
4 their connectivity, and their positions. Phase 3 starts with the optimized mesh (a
5 piecewise linear surface) that is produced in phase 2 and fits an accurate, concise
6 piecewise smooth subdivision surface, again by optimizing an energy function that
7 trades off conciseness and fit to the data. The phase 3 optimization varies the
8 number of vertices in the control mesh, their connectivity, their positions, and the
9 number and locations of sharp features. The automatic detection and recovery of
10 sharp features in the surface is an essential part of this phase.

11 In the present embodiment, processing differs from the approach described
12 in Hoppe et al. in a couple of ways. First, continuous optimization is performed
13 only over vertex positions, since we do not want to alter the connectivity of the
14 control mesh. Additionally, feature constraints are added as well as a smoothing
15 term.

16 In the illustrated example, the fitting process minimizes the functional:

$$18 E(\mathbf{v}) = E_d(\mathbf{v}, \mathbf{p}) + \lambda E_s(\mathbf{v}) + \mu E_c(\mathbf{v})$$

19
20 where \mathbf{v} is a vector of all the vertex positions, and \mathbf{p} is a vector of all the data
21 points from the range scanner. The subscripts on the three terms stand for
22 distance, shape, and constraints. The distance functional E_d measures the sum-
23 squared distance from the range scanner points to the subdivision surface:

$$24 25 E_d(\mathbf{v}, \mathbf{p}) = \sum_{i=1}^{n_p} a_i \|p_i - \Pi(\mathbf{v}, p_i)\|^2$$

1
2 where p_i is the i^{th} range point and $\Pi(\mathbf{v}, p_i)$ is the projection of that point onto the
3 subdivision surface defined by the vertex positions \mathbf{v} . The weight a_i is a Boolean
4 term that causes points for which the scanner's view direction at p_i is not consistent
5 with the surface normal at $\Pi(\mathbf{v}, p_i)$ to be ignored. Additionally, points are rejected
6 that are farther than a certain distance from the surface:

$$7 \quad a_i = \begin{cases} 1 & \text{if } \langle s(p_i), n(\Pi(\mathbf{v}, p_i)) \rangle > 0 \text{ and } \|p_i - \Pi(\mathbf{v}, p_i)\| < d_0 \\ 0 & \text{otherwise} \end{cases}$$

8
9

10 where $s(p)$ is the direction toward the scanner's viewpoint at point p and $n(x)$ is the
11 outward-facing surface normal at point x .

12 The smoothness functional E_s encourages the control mesh to be locally
13 planar. It measures the distance from each vertex to the average of the
14 neighboring vertices:

$$15 \quad E_s(\mathbf{v}) = \sum_{j=1}^{n_v} \left\| \mathbf{v}_j - \frac{1}{\deg(\mathbf{v}_j)} \sum_{i=1}^{\deg(\mathbf{v}_j)} \mathbf{v}_i \right\|^2$$

16
17

18 The vertices \mathbf{v}_i are the neighbors of \mathbf{v}_j .

19 The constraint functional E_c is simply the sum-squared distance from a set
20 of constrained vertices to a set of corresponding target positions:

$$21 \quad E_c(\mathbf{v}) = \sum_{i=1}^{n_c} \|A_c \mathbf{v} - \mathbf{d}_i\|^2$$

22

23 where A_j is the linear function that defines the limit position of the j^{th} vertex
24 in terms of the control mesh, so the limit position of vertex c_i is attached to the 3D
25 point d_i . The constraints could instead be enforced rigidly by a linear

1 reparameterization of the optimization variables, but it has been found that the
2 soft-constraint approach helps guide the iteration smoothly to a desirable local
3 minimum. The constraints are chosen by the user to match the facial features of
4 the generic mesh to the corresponding features on the particular face being fit. In
5 the present example, approximately 25 to 30 constraints are used, concentrating on
6 the eyes, nose, and mouth. Fig. 9 shows the constraints on the subdivision control
7 mesh at 900 and their corresponding points on a face model.

8 Minimizing $E(v)$ is a nonlinear least-squares problem, because Π and a_i are
9 not linear functions of v . However, such can be made a linear problem by holding
10 a_i constant and approximating $\Pi(v, p_i)$ by a fixed linear combination of control
11 vertices. The fitting process therefore proceeds as a sequence of linear least-
12 squares problems with the a_i and the projections of the p_i onto the surface being
13 recomputed before each iteration. The subdivision limit surface is approximated
14 for these computations by the mesh at a particular level of subdivision. Fitting a
15 face takes a small number of iterations (fewer than 30), and the constraints are
16 updated according to a simple schedule as the iteration progresses, beginning with
17 a high λ and low μ to guide the optimization to a very smooth approximation of
18 the face, and progressing to a low λ and high μ so that the final solution fits the
19 data and the constraints closely. The computation time in practice is dominated by
20 computing $\Pi(v, p_i)$.

21 To produce the mesh for rendering, the surface is subdivided to the desired
22 level, producing a mesh that smoothly approximates the face shape. A
23 displacement is then computed for each vertex by intersecting the line normal to
24 the surface at that vertex with the triangulated surface defined by the original scan
25 as described in Lee et al., *Displaced Subdivision Surfaces*, (SIGGRAPH '00

1 Proceedings) July 2000. The resulting surface reproduces all the salient features
2 of the original scan in a mesh that has somewhat fewer triangles, since the base
3 mesh has more triangles in the more important regions of the face. The
4 subdivision-based representation also provides a parameterization of the surface
5 and a built-in set of multiresolution basis functions defined in that
6 parameterization and, because of the feature constraints used in the fitting, creates
7 a natural correspondence across all faces that are fit using this method. This
8 structure is useful in many ways in facial animation.

9 Fig. 10 is a flow diagram that describes steps in a method for building a
10 face model in accordance with this described embodiment. The method can be
11 implemented in any suitable hardware, software, firmware or combination thereof.
12 In the present example, the method is implemented in software.

13 Step 1000 measures 3D data for one or more faces to provide
14 corresponding face models. In the above example, the 3D data was generated
15 through the use of a laser range scan of the faces. It will be appreciated that any
16 suitable method of providing the 3D data can be used. Step 1002 defines a generic
17 face model that is to be used to fit to the one or more face models. It will be
18 appreciated that the generic face model can advantageously be utilized to fit to
19 many different faces. Accordingly, this constitutes an improvement over past
20 methods in which this was not done. In the example described above, the generic
21 face model comprises a mesh structure in the form of a coarse triangle mesh. The
22 triangle mesh defines subdivision surfaces that closely approximate the geometry
23 of the face. In the illustrated example, a single base mesh is used to define the
24 subdivision surfaces for all of the face models. Step 1004 selects specific points
25 or constraints on the generic face model. These specific points or constraints are

1 mapped directly to corresponding points that are marked on the face model. The
2 mapping of these specific points takes place in the same manner for each of the
3 many different possible face models. Step 1006 fits the generic face model to the
4 one or more face models. This step is implemented by manipulating only the
5 positions of the vertices to adapt to the shape of each different face. During the
6 fitting process continuous optimization is performed only over the vertex positions
7 so that the connectivity of the mesh is not altered. In addition, the fitting process
8 involves mapping the specific points or constraints directly to the face model. In
9 addition, a smoothing term is added and minimized so that the control mesh is
10 encouraged to be locally planar.

11

12 **Adding Eyes**

13 The displaced subdivision surface just described represents the shape of the
14 facial skin surface quite well. There are, however, several other features that are
15 desirable for a realistic face. The most important of these is the eyes. Since the
16 laser range scanner does not capture suitable information about the eyes, the mesh
17 is augmented for rendering by adding separately modeled eyes. Unlike the rest of
18 the face model, the eyes and their motions are not measured from a specific
19 person, so they do not necessarily reproduce the appearance of the real eyes.
20 However, their presence and motion is critical to the overall appearance of the face
21 model.

22 Any suitable eye model can be used to model the eyes. In the illustrated
23 example, a commercial modeling package was used to build a model consisting of
24 two parts. The first part is a model of the eyeball, and the second part is a model
25 of the skin surface around the eye, including the eyelids, orbit, and a portion of the

1 surrounding face (this second part will be called the "orbit surface"). In order for
2 the eye to become part of the overall face model, the orbit surface must be made to
3 fit the individual face being modeled and the two surfaces must be stitched
4 together. This is done in two steps: first the two meshes are warped according to a
5 weighting function defined on the orbit surface, so that the face and orbit are
6 coincident where they overlap. Then the two surfaces are cut with a pair of
7 concentric ellipsoids and stitched together into a single mesh.

8 Note that one of the advantageous features of the embodiments described
9 above is that they provide a structure or framework that can be used to transform
10 the expressions of one person into expressions of another person. Because the fit
11 of the generic face model to each individual face is constrained so that any given
12 part of the generic model always maps to the same feature on every person's
13 face—for example, the left corner of the mouth in the generic model always maps
14 to the left corner of the mouth on any person's face—the set of fitted face models
15 provides a means for determining the point on any face that corresponds to a
16 particular point on a particular face. For example, suppose the motion of the left
17 corner of the mouth on person A's face has been measured. We can use the fit of
18 the generic model to face A to determine which point of the generic model
19 corresponds to that measured point, and then we can use the fit of the generic
20 model to face B to determine which point on B's face corresponds to the computed
21 point on the generic model and therefore also to the measured point on face A.
22 This information is essential to transforming motion from one face to another
23 because we have to know which parts of the new face need to be moved to
24 reproduce the motions we measured from a set of points on the measured face.

1 **Moving the Face**

2 The motions of the face are specified by the time-varying 3D positions of a
3 set of sample points on the face surface. When the face is controlled by motion-
4 capture data these points are the markers on the face that are tracked by the motion
5 capture system. The motions of these points are used to control the face surface
6 by way of a set of control points that smoothly influence regions of the surface.
7 Capturing facial motion data can be done in any suitable way, as will be apparent
8 to those of skill in the art. In one specific example, facial motion was captured
9 using the technique described in Guenter et al., *Making Faces*, Proceedings of
10 SIGGRAPH 1998, pages 55-67, 1998.

11 **Mesh Deformation**

12 The face is animated by displacing each vertex w_i of the triangle mesh from
13 its rest position according to a linear combination of the displacements of a set of
14 control points q_j . These control points correspond one-to-one with the sample
15 points p_j that describe the motion. The influence of each control point on the
16 vertices falls off with distance from the corresponding sample point, and where
17 multiple control points influence a vertex, their weights are normalized to sum to
18 1.

20
$$\Delta w_i = \frac{1}{\beta_i} \sum_j \alpha_{ij} \Delta q_j \quad ; \alpha_{ij} = h \left(\frac{\|w_i - p_j\|}{r} \right)$$

21 where $\beta_i = \sum_k \alpha_{ik}$ if vertex i is influenced by multiple control points and 1
22 otherwise. These weights are computed once, using the rest positions of the
23 sample points and face mesh, so that moving the mesh for each frame is just a
24
25

1 sparse matrix multiplication. For the weighting function, the following was used:

2
$$h(x) = \frac{1}{2} + \frac{1}{2}\cos(\pi x).$$

3 Two types of exceptions to these weighting rules are made to handle the
4 particulars of animating a face. Vertices and control points near the eyes and
5 mouth are tagged as "above" and "below," and control points that are, for example,
6 above the mouth do not influence the motions of vertices below the mouth. Also,
7 a scalar texture map in the region around the eyes is used to weight the motions so
8 that they taper smoothly to zero at the eyelids. To move the face mesh according
9 to a set of sample points, control point positions must be computed that will
10 deform the surface appropriately. Using the same weighting functions described
11 above, we compute how the sample points move in response to the control points.
12 The result is a linear transformation: $\mathbf{p} = \mathbf{A}\mathbf{q}$. Therefore if at time t we want to
13 achieve the sample positions \mathbf{p}_t , we can use the control positions $\mathbf{q}_t = \mathbf{A}^{-1}\mathbf{p}_t$.
14 However, the matrix \mathbf{A} can be ill-conditioned, so to avoid the undesirable surface
15 shapes that are caused by very large control point motions we compute \mathbf{A}^{-1} using
16 the SVD (Singular Value Decomposition) and clamp the singular values of \mathbf{A}^{-1} at a
17 limit M . In the illustrated example, $M = 1.5$ was used. A standard reference that
18 discusses SVD is Golub and Van Loan, *Matrix Computations*, 3rd edition, Johns
19 Hopkins press, 1996.

20

21 Eye and Head Movement

22 In order to give the face a more lifelike appearance, procedurally generated
23 motion is added to the eyes and separately captured rigid-body motion to the head
24 as a whole. The eyeballs are rotated according to a random sequence of fixation
25 directions, moving smoothly from one to the next. The eyelids are animated by

1 rotating the vertices that define them about an axis through the center of the
2 eyeball, using weights defined on the eyelid mesh to ensure smooth deformations.

3 The rigid-body motion of the head is captured from the physical motion of
4 a person's head by filming that motion while the person is wearing a hat marked
5 with special machine-recognizable targets (the hat is patterned closely on the one
6 used by Marschner et al., *Image-based BRDF measurement including human skin*,
7 Rendering Techniques '99 (Proceedings of the Eurographics Workshop on
8 Rendering), pps. 131-144, June 1998. By tracking these targets in the video
9 sequence, the rigid motion of the head is computed, which is then applied to the
10 head model for rendering. This setup, which requires simply a video camera,
11 provides a convenient way to author head motion by demonstrating the desired
12 actions.

13

14 **Exemplary System and Methods for Modeling Reflectance**

15 Rendering a realistic image of a face requires not just accurate geometry,
16 but also accurate computation of light reflection from the skin. In the illustrated
17 example, a physically-based Monte Carlo ray tracer was used to render the face.
18 Exemplary techniques are described in Cook et al., *Distribution Ray Tracing*,
19 Computer Graphics (SIGGRAPH '84 Proceedings), pps. 165-174, July 1984 and
20 Shirley et al., *Monte Carlo techniques for direct lighting calculations*,
21 Transactions on Graphics, 15(1):1-36, 1996. Doing so allows for the use of
22 arbitrary BRDFs (bi-directional reflectance distribution functions) to correctly
23 simulate the appearance of the skin, which is not well approximated by simple
24 shading models. In addition, extended light sources are used, which, in rendering
25 as in portrait photography, are needed to achieve a pleasing image. Two important

1 deviations from physical light transport are made for the sake of computational
2 efficiency: diffuse interreflection is disregarded, and the eyes are illuminated
3 through the cornea without refraction.

4 In the illustrated example, a reflectance model for the skin is based on
5 measurements of actual human faces. Exemplary techniques are described in
6 Marschner et al., *Image based BRDF measurement including human skin*,
7 *Rendering Techniques '99* (Proceedings of the Eurographics Workshop on
8 *Rendering*), pps. 131-144, June 1999. The measurements describe the average
9 BRDFs of several subjects' foreheads and include fitted parameters for the BRDF
10 model described in Lafourte et al., *Non-linear approximation of reflectance*
11 *functions*, *Computer Graphics (SIGGRAPH '97 Proceedings)*, pps. 117-126,
12 August 1997. Accordingly, the measurements provide an excellent starting point
13 for rendering a realistic face. However, the measurements need to be augmented
14 to include some of the spatial variation observed in actual faces. This is achieved
15 by starting with the fit to the measured BRDF of one subject whose skin is similar
16 to the skin of the face we rendered and dividing it into diffuse and specular
17 components. A texture map is then introduced to modulate each.

18 The texture map for the diffuse component, or the "albedo map", modulates
19 the diffuse reflectance according to measurements taken from the subjects' actual
20 faces as described below. The specular component is modulated by a scalar
21 texture map to remove specularity from areas (such as eyebrows and hair) that
22 should not be rendered with skin reflectance and to reduce specularity on the
23 lower part of the face to approximate the characteristics of facial skin. The result
24 is a spatially varying BRDF that is described at each point by a sum of the
25 generalized cosine lobes of Lafourte et al., *Non-linear approximation of*

1 reflectance functions, Computer Graphics (SIGGRAPH '97 Proceedings), pps.
2 117-126, August 1997.

3

4 **Constructing the Albedo Map**

5 In the illustrated and described embodiment, the albedo map, which must
6 describe the spatially varying reflectance due to diffuse reflection, was measured
7 using a sequence of digital photographs of the face taken under controlled
8 illumination.

9 Fig. 11 shows an exemplary system that was utilized to capture the digital
10 photographs or images. In the illustrated system, a digital camera 1100 is
11 provided and includes multiple light sources, exemplary ones of which are shown
12 at 1102, 1104. Polarizing filters in the form of perpendicular polarizers 1106,
13 1108, and 1110 are provided and cover the light sources and the camera lens so
14 that the specular reflections are suppressed, thereby leaving only the diffuse
15 component in the images. In the example, a subject wears a hat 1112 printed with
16 machine-recognizable targets to track head pose. Camera 1100 stays stationary
17 while the subject rotates. The only illumination comes from the light sources
18 1102, 1104 at measured locations near the camera. A black backdrop is used to
19 reduce indirect reflections from spilled light.

20 Since the camera and light source locations are known, standard ray tracing
21 techniques can be used to compute the surface normal, the irradiance, the viewing
22 direction, and the corresponding coordinates in texture space for each pixel in each
23 image. Under the assumption that ideal Lambertian reflection is being observed,
24 the Lambertian reflectance can be computed for a particular point in texture space
25 from this information. This computation is repeated for every pixel in one

photograph which essentially amounts to projecting the image into texture space and dividing by the computed irradiance due to the light sources to obtain a map of the diffuse reflectance across the surface. Consider Fig. 12 in which two photographs are shown projected into texture space and corrected for lighting. In practice the projection is carried out by reverse mapping, with the outer loop iterating through all the pixels in the texture map, and stochastic supersampling is used to average over the area in the image that projects to a particular texture pixel.

The albedo map from a single photograph only covers part of the surface, and the results are best at less grazing angles. Accordingly a weighted average of all the individual maps is computed to create a single albedo map for the entire face. The weighting function, a visual example of which is given in Fig. 13, should be selected so that higher weights are given to pixels that are viewed and/or illuminated from directions nearly normal to the surface, and should drop to zero well before either viewing or illumination becomes extremely grazing. In the illustrated example, the following function was used $(\cos \theta_i \cos \theta_e - c)^p$, with $c = 0.2$ and $p = 4$.

Before computing the albedo for a particular texture pixel, we verify that the pixel is visible and suitably illuminated. Multiple rays are traced from points on the pixel to points on the light source and to the camera point, and the pixel is marked as having zero, partial, or full visibility and illumination. It is prudent to err on the large side when estimating the size of the light source. Only albedos for pixels that are fully visible, fully illuminated by at least one light source, and not partially illuminated by any light source are computed. This ensures that partially

1 occluded pixels and pixels that are in full-shadow or penumbra regions are not
2 used.

3 Some calibration is required to make these measurements meaningful. The
4 camera's transfer curve was calibrated using the method described in Debevec et
5 al., *Recovering high dynamic range radiance maps from photographs*, Computer
6 Graphics (SIGGRAPH '97 Proceedings), pps. 369-378, August 1997. The
7 light/camera system's flat-field response was calibrated using a photograph of a
8 large white card. The lens's focal length and distortion were calibrated using the
9 technique described in Zhang, *A flexible new technique for camera calibration*,
10 Technical Report MSR-TR-98-71, Microsoft Research, 1998. The absolute scale
11 factor was set using a reference sample of known reflectance. When image-to-
12 image variation in light source intensity was a consideration, control was provided
13 by including the reference sample in every image.

14 The texture maps that result from this process do a good job of
15 automatically capturing the detailed variation in color across the face. In a few
16 areas, however, the system cannot compute a reasonable result. Additionally, the
17 strap used to hold the calibration hat in place is visible. These problems are
18 removed by using an image editing tool and filling in blank areas with nearby
19 texture or with uniform color.

20 Figs. 14 and 15 show the raw and edited albedo maps for comparison. The
21 areas where the albedo map does not provide reasonable results can be seen where
22 the surface is not observed well enough (e. g., under the chin) or is too intricately
23 shaped to be correctly scanned and registered with the images (e.g the ears).
24 Neither of these types of areas requires the texture from the albedo map for
25 realistic appearance—the first because they are not prominently visible and the

1 second because the geometry provides visual detail—so this editing has relatively
2 little effect on the appearance of the final renderings.

3 Fig. 16 shows several different aspects of the face model, using still frames
4 from the accompanying video. In the first row, the face is shown from several
5 angles to demonstrate that the albedo map and measured BRDF realistically
6 capture the distinctive appearance of the skin and its color variation over the entire
7 face, viewed from any angle. The second row shows the effects of rim and side
8 lighting, including strong specular reflections at grazing angles. Note that the light
9 source has the same intensity and is at the same distance from the face for all three
10 images in this row. The directional variation in the reflectance leads to the
11 familiar lighting effects seen in the renderings. In the third row, expression
12 deformations are applied to the face to demonstrate that the face still looks natural
13 under normal expression movement.

14 Fig. 17 is a flow diagram that describes steps in a method for creating an
15 albedo map in accordance with the described embodiment. The method can be
16 implemented in any suitable hardware, software, firmware or combination thereof.
17 In the described embodiment, the method is implemented in software in
18 connection with a system such as the one shown and described in Fig. 11.

19 Step 1700 provides one or more polarized light sources that can be used to
20 illuminate a subject. Exemplary light sources are described above. In the
21 described embodiment, the light sources are selected so that the specular
22 component of the subject's facial reflectance is suppressed or eliminated. Step
23 1702 illuminates the subject's face with the light sources. Step 1704 rotates the
24 subject while a series of digital photographs or images are taken. Step 1706
25 computes surface normals, irradiance, viewing direction and coordinates in texture

space for each pixel in the texture map. The computations can be done using known algorithms. Step 1708 computes the Lambertian reflectance for a particular pixel in the texture space for the image. This provides an albedo for the pixel. Step 1710 determines whether there are any additional pixels in the albedo map. If there are, step 1712 gets the next pixel and returns to step 1708. If there are no additional pixels in the albedo map, step 1714 ascertains whether there are any additional digital images. If there are additional digital images, step 1716 gets the next digital image and returns to step 1706. If there are no additional digital images, then step 1718 computes a weighted average of the individual albedo maps for each image to create a single albedo map for the entire face. One specific example of how this weighted average processing takes place is given above and described in Marschner, *Inverse Rendering for Computer Graphics*, PhD thesis, Cornell University, August 1998.

Fig. 18 is a flow diagram that describes steps in a method for computing an albedo for a single pixel. This method can be implemented in any suitable hardware, software, firmware or combination thereof. In the described embodiment, the method is implemented in software. Step 1800 determines, for a given pixel, whether the pixel is fully visible. If the pixel is not fully visible, then an albedo for the pixel is not computed (step 1804). If the pixel is fully visible, step 1802 determines whether the pixel is fully illuminated by at least one light source. If the pixel is not fully illuminated by at least one light source, then an albedo for the pixel is not computed (step 1804). If the pixel is fully illuminated by at least one light source, then step 1806 determines whether the pixel is partially illuminated by any light source. If so, then an albedo is not computed for the pixel. If the pixel is not partially illuminated by any light source, then step

1 1808 computes an albedo and a weight for the pixel. The weights are later used in
2 averaging together individual maps. Hence, as discussed above, albedos are
3 computed only for pixels that are fully visible, fully illuminated by at least one
4 light source, and not partially illuminated by any light source. This ensures that
5 partially occluded pixels and pixels that are in full-shadow or penumbra are not
6 used.

7

8 Conclusion

9 The embodiments described above provide systems and methods that
10 address the challenge of modeling and rendering faces to the high standard of
11 realism that must be met before an image as familiar as a human face can appear
12 believable. The philosophy of the approach is to use measurements whenever
13 possible so that the face model actually resembles a real face. The geometry of the
14 face is represented by a displacement-mapped subdivision surface that has
15 consistent connectivity and correspondence across different faces. The reflectance
16 comes from previous BRDF measurements of human skin together with new
17 measurements that combine several views into a single illumination-corrected
18 texture map for diffuse reflectance. The motion comes from previously described
19 motion capture technique and is applied to the face model using an improved
20 deformation method that produces motions suitable for shaded surfaces. The
21 realism of the renderings is greatly enhanced by using the geometry, motion, and
22 reflectance of real faces in a physically-based renderer.

23 Although the invention has been described in language specific to structural
24 features and/or methodological steps, it is to be understood that the invention
25 defined in the appended claims is not necessarily limited to the specific features or

1 steps described. Rather, the specific features and steps are disclosed as preferred
2 forms of implementing the claimed invention.

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25